# Architecture trade-offs in a planet-scale queueing system
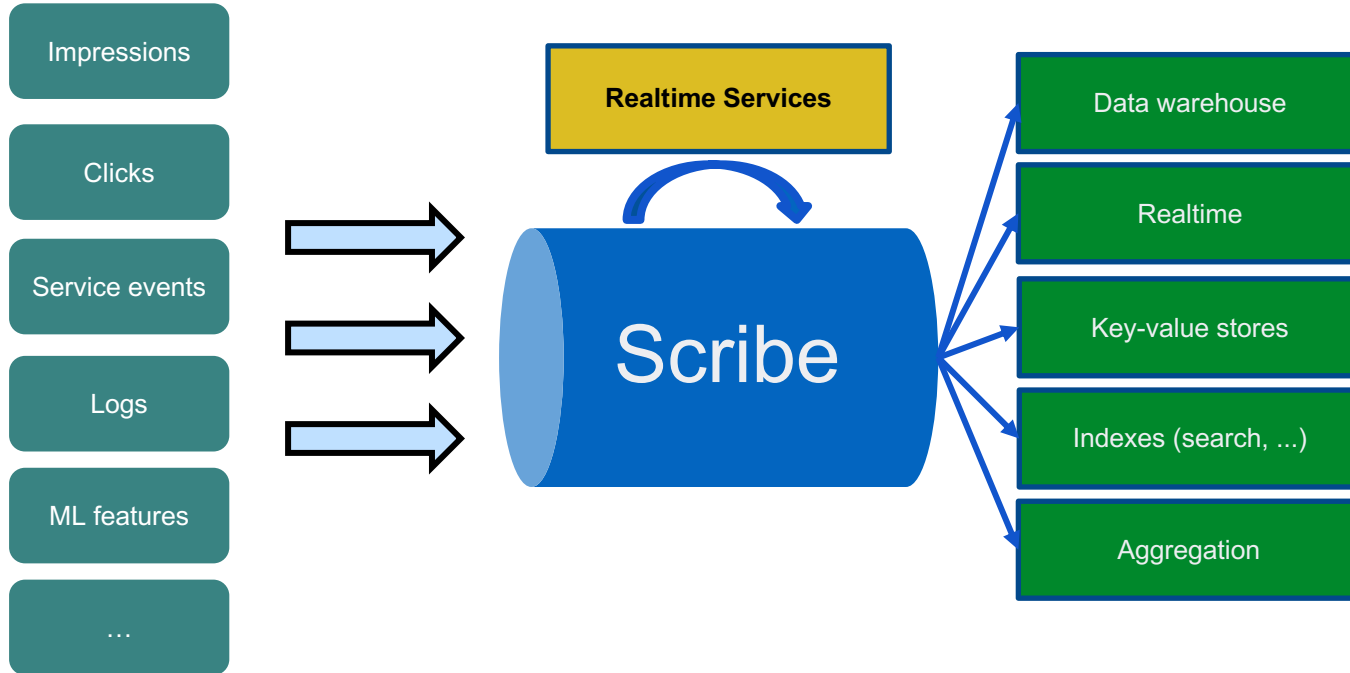
**Manos Karpathiotakis**

facebook
London Office

# Getting data from point A to point B

# Getting data from point A to point B



**Distributed, buffered, multi-tenant pipe**

# Hello World

```
manos at vm4 > scribe_cat testcat hello
manos at vm4 > scribe_cat testcat world
manos at vm4 > ▐
```

```
manos at vm5 > scribe_cat testcat foo
manos at vm5 > scribe_cat testcat bar
manos at vm5 > ▐
```
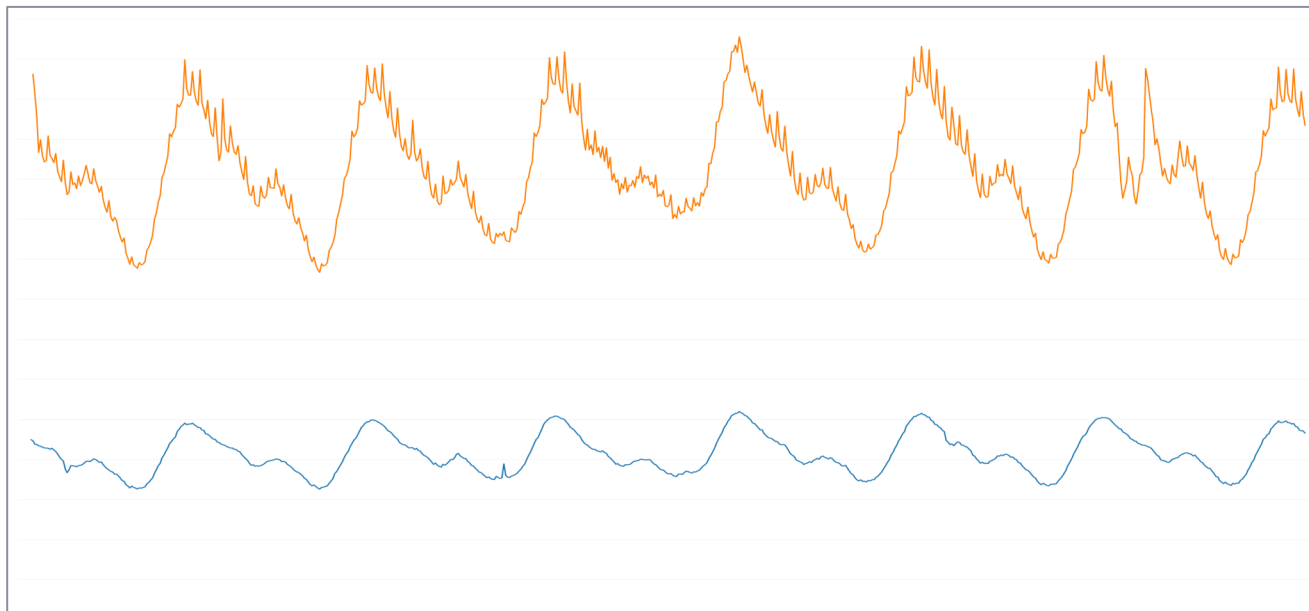
```
manos at vm6 in ~ > ptail -f testcat
hello
foo
bar
world
```

**Logical stream abstraction**

# Scale

Read: 7TB/s

Write: 2.5TB/s



Millions of machines

Hundreds of thousands of categories/topics

**Scale does not come for free**

4

Retention in
the days

Available on
every machine

Latency seconds
to minutes

Slightly lossy,
3 9s to 5 9s

Rough ordering
guarantees

Retention in the days

Available on every machine

Latency seconds to minutes

Slightly lossy, 3 9s to 5 9s
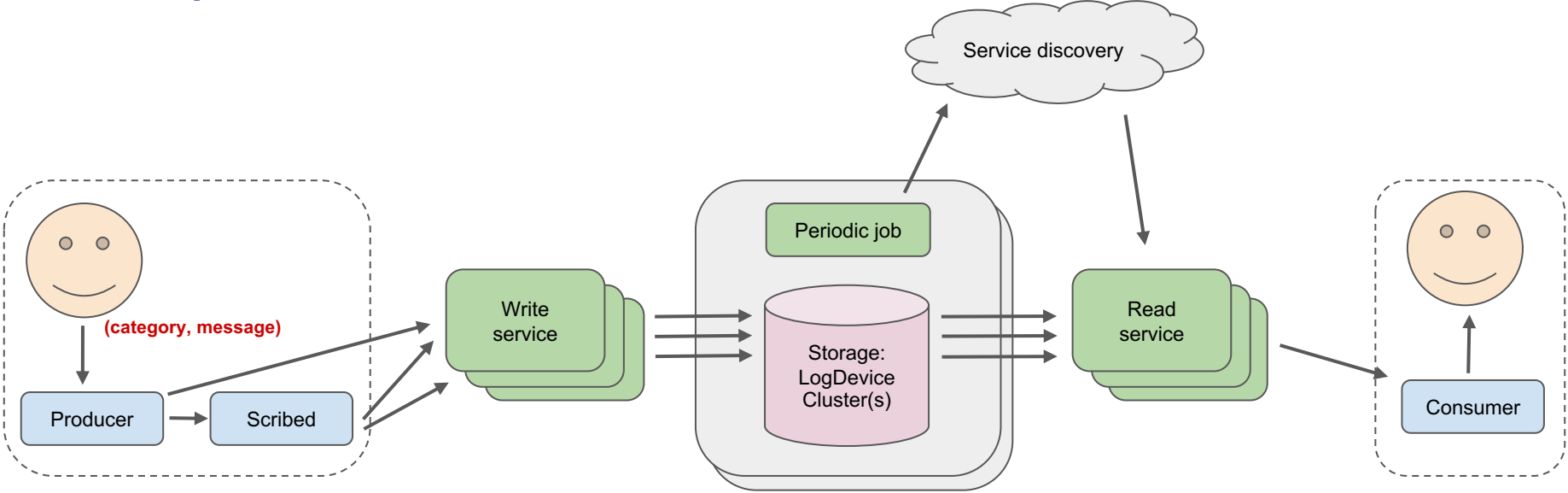
Rough ordering guarantees

**This talk**

# Data Completeness

# Customers willing to lose data?!

- High volumes generated in real time ("nowhere to park the data")

- Lossy or sampled upstream ("not making things worse")

- Statistical in nature where small losses not critical (ML use cases)

- Data freshness imperative ("stale data is useless data")

**Multiple customers are unwilling to "pay" for completeness**

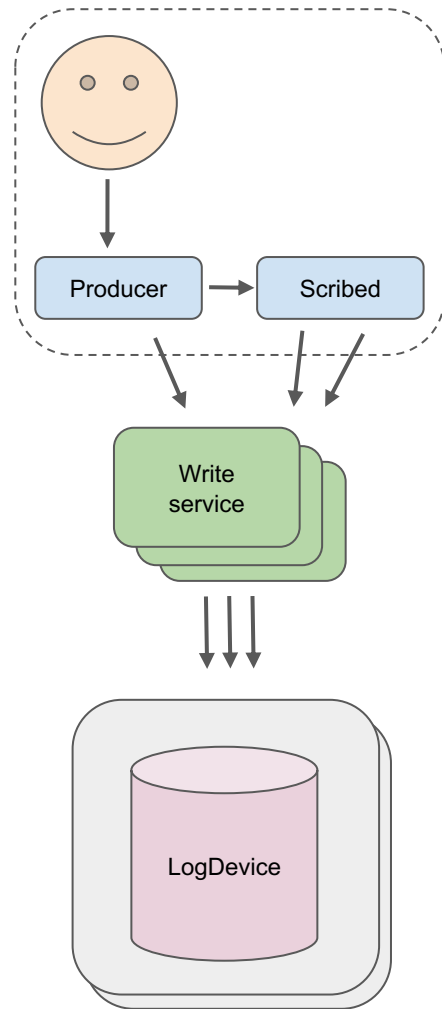# Data plane



**(Big) Aggregation Tree**

# Write path

LogDevice is durable storage

But on the way, single copy in memory

When to acknowledge producer operations?

1. Once processed in the producer?

2. Once processed in the write service?

3. Once stored in LogDevice?
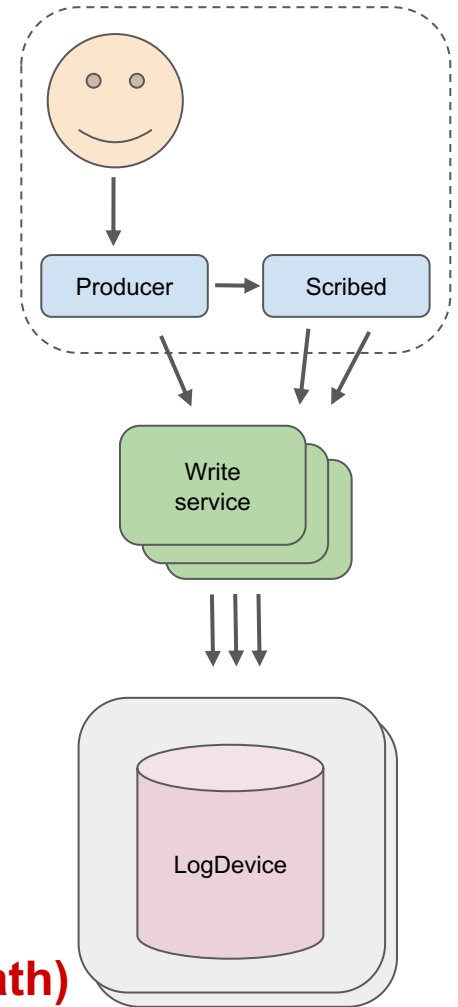
**Offer customers multiple flavors**

# High durability flavor

Acknowledgement once stored in LogDevice

Increased duplication

Less aggressive batching => Lower throughput

Lower latency

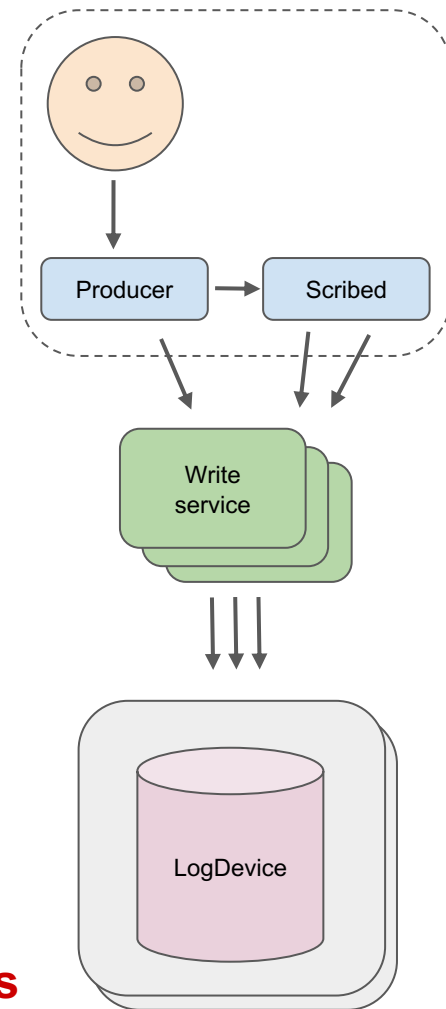**At least once semantics (…on write path)**

# High throughput flavor

Producer (**optionally**) acknowledges upon receival

Accept small amount of data loss

Heavy batching provides high scalability

"Approximately once" semantics



**"Approximately once" semantics**
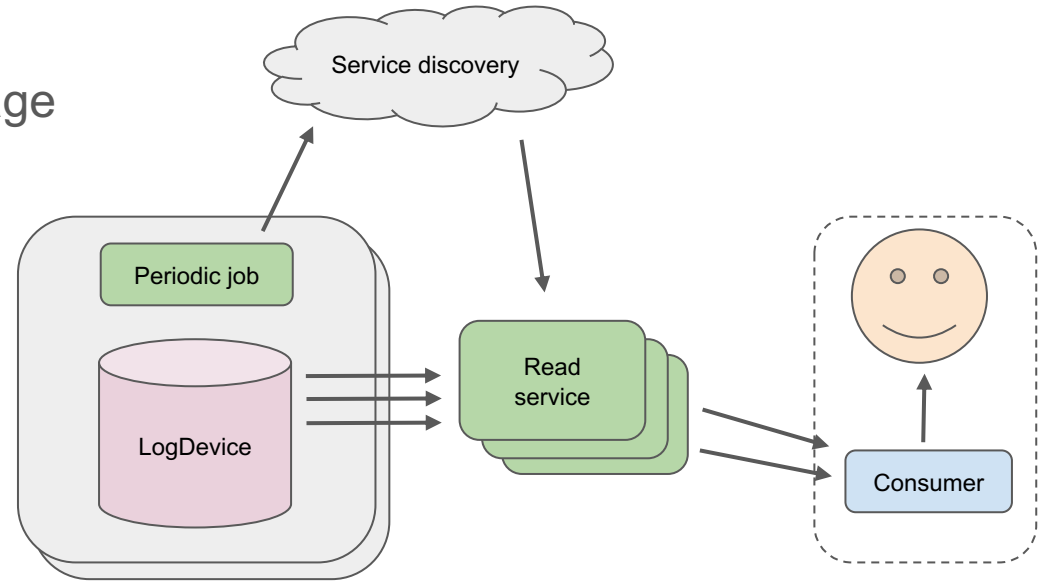
# Read path

Single logical copy for each message

What if cluster unavailable:

- Data unavailable

Options

- Accept loss and carry on

- Wait

- Abort



**"Clean" layering minimizes complexity yet is prone to data loss**
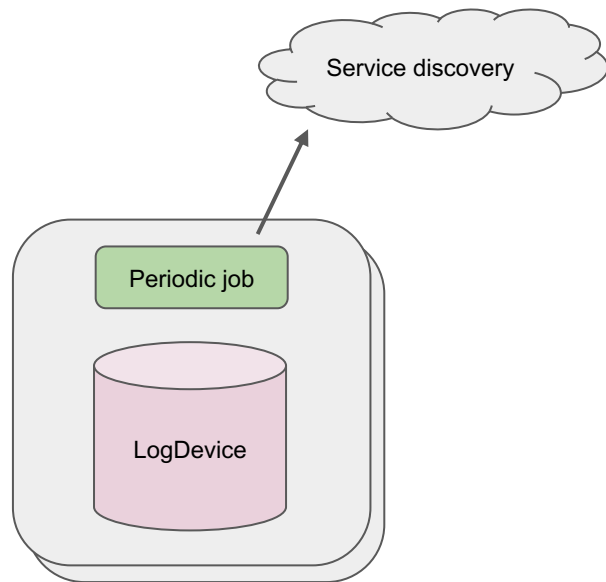
# (Rough) Ordering

# Storage

Each LogDevice cluster has a set of configured logs

- Each log holds data for only one category
- One category can have multiple logs
- Logs have a maximum throughput

Periodic job responsibilities

- Publish which categories have data
- Split logs if they get hot



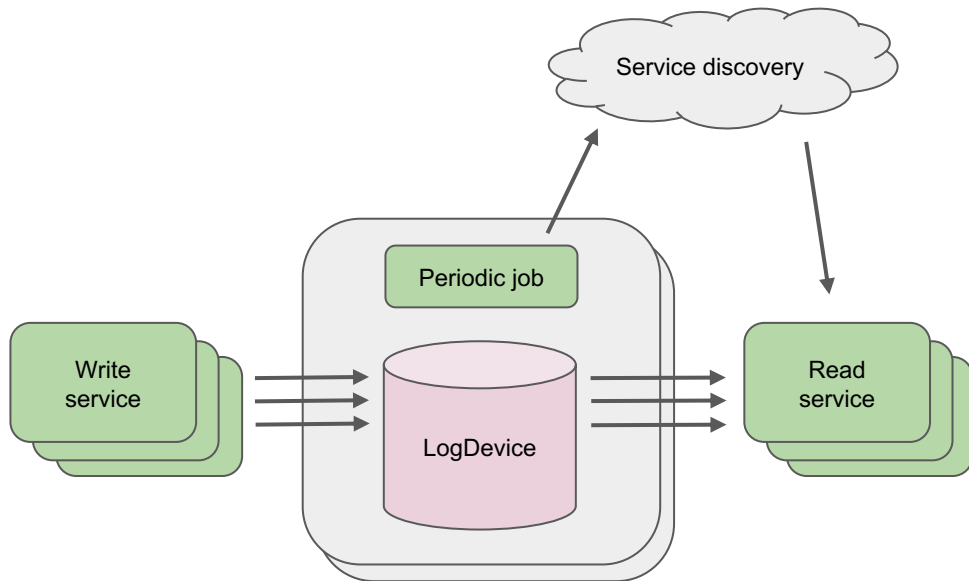Service discovery

Periodic job

LogDevice

# Storage

When Write Service sees (cat, msg)

- Pick a cluster
- Pick any log for cat
- Append msg

When read service sees (cat, time)

- Lookup clusters with data for cat
- Lookup time in all relevant logs
- Merge streams for logs into single output

# Properties

Traffic to a single category scales horizontally

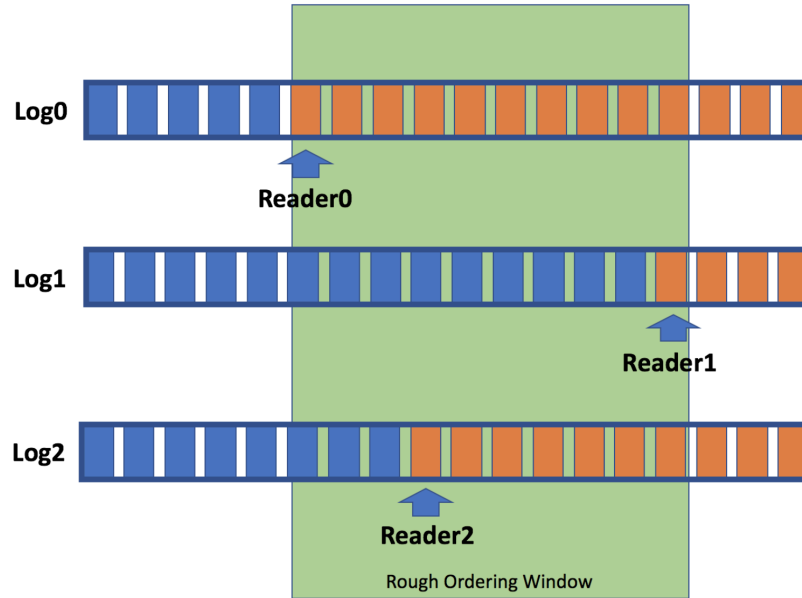Consecutive writes can end up in entirely different clusters

Top priorities

- Scalability
- (Write) availability

Lost in the process

- Ordering guarantees
- Repeatable reads

# Rough ordering



**Reduce blast radius of stragglers**

**Relaxed semantics in favor of read availability**

# There is no free lunch

Inherent trade-off between **scale**, operational **complexity**, and **semantics**

- Often, the semantics are held constant ("no loss", "strict ordering")

- In Scribe's case, scale is imposed by company growth

- Relaxing semantics as a tool to manage complexity

- Users can still build (more) reliable apps over Scribe (at an extra cost)

**When relaxing semantics, let users decide**

# Further information

[2019] Facebook eng blog post

engineering.fb.com/data-infrastructure/scribe

[2019] Tech talk Systems@Scale NYC

facebook.com/atscaleevents/videos/509450066277552

[2016] Realtime Data Processing at Facebook

research.fb.com/publications/realtime-data-processing-at-facebook

# Questions?

# Background

Scribe has been around for 10+ years

Initial purpose was to batch and store logs

Purpose evolved a lot over the years

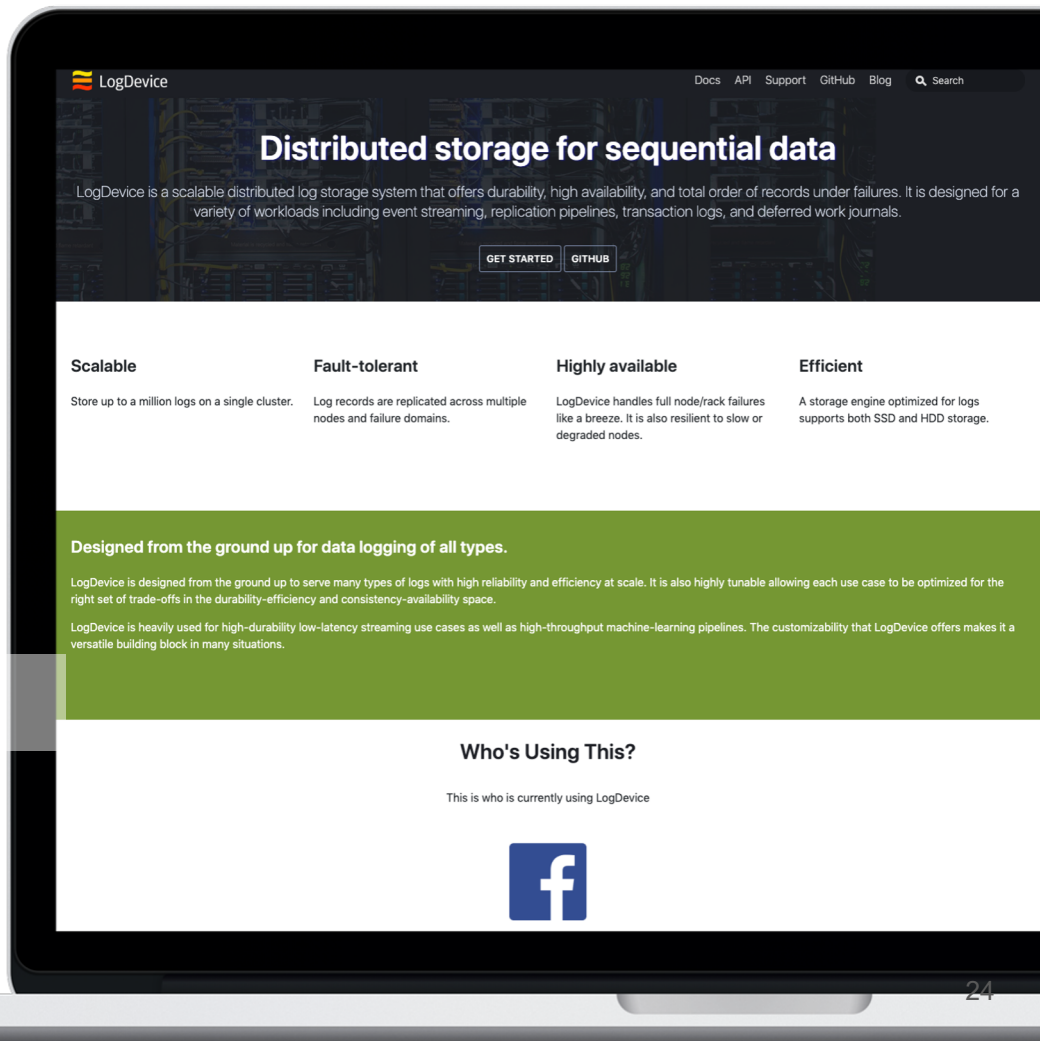Rearchitected multiple times to cope with scale
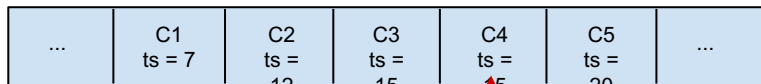
# LogDevice

Distributed storage

Log as a primitive

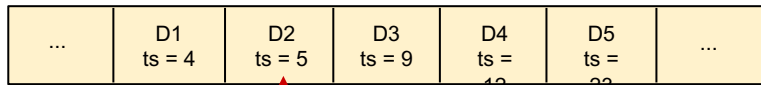Ideal for streams of data

logdevice.io

# Rough ordering

Log C, in cluster X

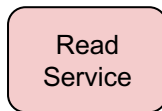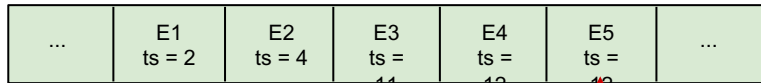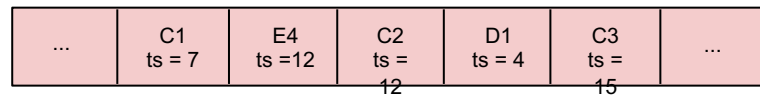| ... | C1 ts = 7 | C2 ts = 12 | C3 ts = 15 | C4 ts = 15 | C5 ts = 20 | ... |
|---|---|---|---|---|---|---|

**Window = 10**

Log D, in cluster X

| ... | D1 ts = 4 | D2 ts = 5 | D3 ts = 9 | D4 ts = 12 | D5 ts = 23 | ... |
|---|---|---|---|---|---|---|

Log E, in cluster Y

| ... | E1 ts = 2 | E2 ts = 4 | E3 ts = 11 | E4 ts = 12 | E5 ts = 12 | ... |
|---|---|---|---|---|---|---|

Read Service

Output

| ... | C1 ts = 7 | E4 ts =12 | C2 ts = 12 | D1 ts = 4 | C3 ts = 15 | ... |
|---|---|---|---|---|---|---|

**Reduce blast radius of stragglers**

**Relaxed semantics in favor of read availability**